

CS 433: Computer Architecture – Fall 2021

Homework 5

Total Points: Undergraduates (44 points), Graduates (52 points)

Undergraduate students should only solve the first 4 problems.

Graduate students should solve all problems.

Due Date: November 11, 2021 at 10:00 pm CT

(See course information slides for more details)

Directions:

- **All students must write and sign the following statement at the end of their homework submission. "I have read the honor code for this class in the course information handout and have done this homework in conformance with that code. I understand fully the penalty for violating the honor code policies for this class." No credit will be given for a submission that does not contain this signed statement.**
- **On top of the first page of your homework solution, please write your name and NETID, your partner's name and NETID, and whether you are an undergrad or grad student.**
- **Name your homework solution file as *firstname_lastname_hw5.pdf***
- **Please show all work that you used to arrive at your answer. Answers without justification will not receive credit. Errors in numerical calculations will not be penalized. Cascading errors will usually be penalized only once.**
- **See course information slides for more details.**

Problem 1 [5 Points]

A 4 entry victim cache for a 4KB direct mapped cache removes 80% of the conflict misses in a program. Without the victim cache, the miss rate is 0.064 (6.4%) and 67% of these misses are conflict misses. What is the percentage improvement in the AMAT (average memory access time) due to the victim cache?

Assume a hit in the main (4KB) cache takes 1 cycle. For a miss in the main cache that hits in the victim cache, assume an additional penalty of 1 cycle to access the victim cache. For a miss in both the main and victim caches, assume a further penalty of 48 cycles to get the data from memory. Assume a simple, single-issue, 5-stage pipeline, in-order processor that blocks on every read and write until it completes.

Problem 2 [12 points]

You are building a computer system around a processor with in-order execution that runs at 1 GHz and has a CPI of 1, excluding memory accesses. The only instructions that read or write data from/to memory are loads (20% of all instructions) and stores (5% of all instructions).

The memory system for this computer has a split L1 cache. Both the I-cache and the D-cache hold 32 KB each. The I-cache has a 2% miss rate and 64 byte blocks, and the D-cache is a write-through, no-write-allocate cache with a 5% miss rate and 64 byte blocks. The hit time for both the I-cache and the D-cache is 1 ns. The L1 cache has a write buffer. 95% of writes to L1 find a free entry in the write buffer immediately. The other 5% of the writes have to wait until an entry frees up in the write buffer (assume that such writes arrive just as the write buffer initiates a request to L2 to free up its entry and the entry is not freed up until the L2 is done with the request). The processor is stalled on a write until a free write buffer entry is available.

The L2 cache is a unified write-back, write-allocate cache with a total size of 512 KB and a block size of 64-bytes. The hit time of the L2 cache is 15ns for both read hits and write hits. Tag comparison for hit/miss is included in the 15ns in all cases, do not add hit time to miss time on a miss. The local hit rate of the L2 cache is 80%. Also, 50% of all L2 cache blocks replaced are dirty. The 64-bit wide main memory has an access latency of 20ns (including the time for the request to reach from the L2 cache to the main memory), after which any number of bus words may be transferred at the rate of one bus word (64-bit) per bus cycle on the 64-bit wide 100 MHz main memory bus. Assume inclusion between the L1 and L2 caches and assume there is no write-back buffer at the L2 cache. Assume a write-back takes the same amount of time as an L2 read miss of the same size.

Assume all caches in the system are blocking; i.e., they can handle only one memory access (load, store, or writeback) at a time. When calculating the miss penalty for a load or store for a writeback cache, the time for any needed writebacks should be included in the miss penalty.

While calculating any time values (such as hit time, miss penalty, AMAT), please use ns (nanoseconds) as the unit of time. For miss rates below, give the **local miss rate** for that cache. By **miss penalty_{L2}**, we mean the time from the miss request issued by the L2 cache up to the time the data comes back to the L2 cache from main memory.

Part A [7 points]

Computing the AMAT (average memory access time) for *instruction accesses*.

- i. Give the values of the following terms for instruction accesses: **hit time_{L1}**, **miss rate_{L1}**, **hit time_{L2}**, **miss rate_{L2}**. [1 point]
- ii. Give the formula for calculating **miss penalty_{L2}**, and compute the value of miss penalty L2. [4 points]

iii. Give the formula for calculating the AMAT for this system using the five terms whose values you computed above and any other values you need. [1 point]

iv. Plug in the values into the AMAT formula above, and compute a numerical value for AMAT for instruction accesses. [1 point]

Part B [2 points]

Computing the AMAT for *data reads*.

i. Give the value of **miss rate**_{L1} for data reads. [1 point]

ii. Calculate the value of the AMAT for data reads using the above value, and other values you need. [1 point]

Part C [3 points]

Computing the AMAT for *data writes*. Assume miss penalty_{L2} for a data write is the same as that computed previously for a data read.

i. Give the value of **write time**_{L2Buff}, the time for a write buffer entry to be written to the L2 cache. [2 points]

ii. Calculate the value of the AMAT for data writes using the above information, and any other values that you need. Only include the time that the processor will be stalled. Hint: There are two cases to be considered here depending upon whether the write buffer is full or not. [1 point]

Problem 3 [13 points]

Consider the following piece of code:

```
register int i, j; /* i, j are in the processor registers */
register float sum1, sum2;
float a[64][64], b[64][64];

for (i = 0; i < 64; i++) {                               /* 1 */
    for (j = 0; j < 64; j++) {                           /* 2 */
        sum1 += a[i][j];                                /* 3 */
    }
    for (j = 0; j < 32; j++) {                           /* 4 */
        sum2 += b[i][2*j];                              /* 5 */
    }
}
```

Assume the following:

- There is a perfect instruction cache; i.e., do not worry about the time for any instruction accesses.
- Both *int* and *float* are of size 4 bytes.
- Only the accesses to the array locations $a[i][j]$ and $b[i][2*j]$ generate loads to the data cache. The rest of the variables are all allocated in registers.
- Assume a fully associative, LRU data cache with 32 lines, where each line has 16 bytes.
- Initially, the data cache is empty.
- To keep things simple, we will assume that statements in the above code are executed sequentially. The time to execute lines (1), (2), and (4) is 4 cycles for each invocation. Lines (3) and (5) take 10 cycles to execute and an additional 40 cycles to wait for the data if there is a data cache miss.
- There is a data prefetch instruction with the format `prefetch(array[index])`. This prefetches the entire block containing the word `array[index]` into the data cache. It takes 1 cycle for the processor to execute this instruction and send it to the data cache. The processor can then go ahead and execute subsequent instructions. If the prefetched data is not in the cache, it takes 40 cycles for the data to get loaded into the cache.
- The arrays **a** and **b** are stored in row major form.
- The arrays **a** and **b** both start at cache line boundaries.

Part A [2 points]

How many cycles does the above code fragment take to execute if we do NOT use prefetching?

Part B [2 points]

Consider inserting prefetch instructions for the two inner loops for the arrays **a** and **b** respectively. Explain why we may want to unroll the loops to insert prefetches. What is the minimum number of times you would need to unroll for each of the two loops for this purpose?

Part C [4 points]

Unroll the inner loops for the number of times identified in part b, and insert the minimum number of software prefetches to minimize execution time. The technique to insert prefetches is analogous to software pipelining. You do not need to worry about startup and cleanup code and do not introduce any new loops.

Part D [2 points]

How many cycles does the code in part (c) take to execute? Calculate the average speedup over the code without prefetching. Assume prefetches are not present in the startup code. Extra time needed by prefetches executing beyond the end of the loop execution time should not be counted.

Part E [3 points]

Is there another technique that can be used to achieve the same objective as loop unrolling in this example, but with fewer instructions? Explain this technique and illustrate its use for the code in part (c).

Problem 4 [14 points]

Way prediction allows an associative cache to provide the hit time of a direct-mapped cache. The MIPS R10000 processor used way prediction to achieve a different goal: reduce the cost of the chip package. The R10000 hardware includes an on-chip L1 cache, on-chip L2 tag comparison circuitry, and an on-chip L2 way prediction table. L2 tag information is brought on chip to detect an L2 hit or miss. The way prediction table contains 8K 1-bit entries, each corresponding to two L2 cache blocks. L2 cache storage is built external to the processor package, is 2-way associative, and may have one of several block sizes.

Part A [2 points]

How can way prediction reduce the number of pins needed on the R10000 package to read L2 tags and data, and what is the impact on performance compared to a package with a full complement of pins to interface to the L2 cache?

Part B [2 points]

How could a 2-associative cache be implemented with the same smaller number of pins but without the way prediction table? What is the performance drawback?

Part C [4 points]

Assume that the R10000 uses most-recently used way prediction. What are reasonable design choices for the cache state update(s) to make when the desired data is in the predicted way, the desired data is in the non-predicted way, and the desired data is not in the L2 cache? Please fill in your answers in the following table.

Cache Access Case	Cache State Change Way Prediction Entry
Desired data is in the predicted way	No change
Desired data is in the non-predicted way	
Desired data is not in the L2 cache	

Part D [2 points]

For a 1024 KB L2 cache with 64-byte blocks and 8-way set associativity, how would the prediction table be organized for this new size? Give your answer in the form of “X entries by Y bits per entry.”

Part E [2 points]

For an 8 MB L2 cache with 128-byte blocks and 2-way set associativity, what would the prediction table organization be? Again, give your answer as “X entries by Y bits per entry.”

Part F [2 points]

What is the difference in the way that the R10000 with only 8K way prediction table entries will support the cache in part d) versus the cache in part e)? Hint: Think about the similarity between a way prediction table and a branch prediction table.

NOTE: ONLY GRADUATE STUDENTS SHOULD SOLVE THIS PROBLEM

Problem 5 [8 points]

Consider a computer with an in-order CPU, and with a data cache block size of 64 bytes (16 words) and a 32-bit wide bus to the memory. The memory takes 10 cycles to supply the first word and 2 cycles per word to supply the rest of the block. The cache is non-blocking, and it can support any number of outstanding misses. The memory can service multiple requests simultaneously if required (techniques to achieve this will be discussed in class).

This cache and memory system implement a “**Requested Word First and Early Restart**” policy, and the bus delivers the block data in “**cyclic order**” starting with the requested word. Cyclic order means that if the requested word is the 5th in a block of size 16 words, then the order in which the words in the block are supplied is 5, 6, 7 ... 16, 1, 2, 3, 4.

Part A [3 points]

Consider the following code fragment, which operates on an integer array A which is block-aligned (that is A[0] is located at the start of a cache block in memory):

```
for (i = 11; i < 100; i += 16) {           /* 1 */
    A[i] *= 2;                             /* 2 */
}
```

Suppose that the cache is big enough so that there are only compulsory misses. Further, statement 1 takes 4 cycles to execute, and statement 2 takes 4 cycles to execute in addition to any miss latency. Assume no overlap in the execution of these statements. Initially, the array A is not present in the cache, so any initial accesses to A cause misses in the cache.

What is the running time of this loop with the “*Requested Word First and Early Restart*” policy?

Part B [3 points]

How many cycles would the above loop take to run in a system with just “*Early Restart*” (i.e. the block is fetched in normal order, but the program is started early at arrival of requested word).

Part C [2 points]

How many cycles would the above loop take to run in a system with the base policy (i.e. normal fetch and restart)?